

CLAIMS

What is claimed is:

1. A system for optimizing server selection for clients from among a plurality of servers in a packet communication network, the system comprising:
 - 5 a plurality of servers for alternatively responding to client requests;
 - a central server that maintains server selection weights, and, based on the weights, provides a candidate server list for responding to a client request, the central server receiving feedback indicating service by individual servers in response to client requests and modifying the server selection weights based on the feedback.
- 10 2. The system according to Claim 1 further comprising a DNS server, the DNS server:
 - receives the client request from the client; and
 - based on the client requests, forwards the client requests to the central server.
- 15 3. The system according to Claim 2, wherein the DNS server interrogates candidate servers in the candidate server list.
4. The system according to Claim 3, wherein the DNS server selects a candidate server based on the interrogation.
- 20 5. The method according to Claim 4, wherein the DNS server:
 - indicates to the selected candidate server that it has been selected to provide service to the requesting client; and
 - returns the address of the selected candidate server to the client.

6. The system according to Claim 3, wherein the DNS server returns to the requesting client the address of the first server to respond to the interrogation.
 7. The system according to Claim 6, wherein the DNS server transmits to the client a redirection packet to inform the selected server of being selected.
- 5 8. The system according to Claim 1, wherein the candidate server list includes extra, randomly selected, candidate servers beyond the candidate servers selected based on the weights.
9. The system according to Claim 8, wherein the extra, randomly selected, candidate servers are a fixed percentage beyond the number of servers selected based on the weights.
- 10 10. The system according to Claim 8, wherein the extra, randomly selected, candidate servers are a fixed number beyond the number of server addresses selected based on the weights.
11. The system according to Claim 1, wherein each candidate server in the candidate server list is unique from each other candidate server in the list.
- 15 12. The system according to Claim 1, wherein the feedback occurs according to at least one of the following: number of times the respective server is selected, duration from last feedback, time of day, or requested event.
13. The system according to Claim 1, wherein the weights are based on bias factors to reduce convergence time, the bias factors including at least one of: a number of times selected, moving average based on a specified amount of time, time of day, time of year, calendar event, or geographical location.
- 20

14. The system according to Claim 1, wherein the probabilities sum to one.
15. The system according to Claim 1, wherein the central server includes vectors of server selection weights for subsets of clients.
16. The system according to Claim 1, wherein the central server includes multiple central servers organized as a distributed system.
5
17. The system according to Claim 1, wherein the client interrogates the candidate servers in the candidate server list to measure at least one of the following: network performance between the client and candidate server, server congestion, or server load.
- 10 18. The system according to Claim 1, wherein the candidates represented in the candidate server list are pseudo-randomly selected based on the weights.
19. A method for optimizing server selection for clients from among a plurality of servers in a packet communication network, the method comprising the steps of:
15 providing a candidate server list by a central server for a client requesting a server address, the candidate server list including server addresses selected based on weights corresponding to the candidate servers;
selecting a preferred server from candidate server list;
feeding back service metrics to the central server corresponding to service provided by the respective servers; and
20 updating the weights based on the service metric.
20. The method according to Claim 19 further including the step of forwarding the client request to the central server by a DNS server.

21. The method according to Claim 20, wherein the step of selecting a preferred server is executed by the DNS server.
22. The method according to Claim 21, further including, by the DNS server, the steps of:
 - 5 informing the selected server that it has been selected to provide service to the requesting client; and
 - returning the address of the selected server to the client.
23. The method according to claim 21, further including, by the DNS server, the step of returning to the requesting client the address of the first server from which the probe is returned.
 - 10
24. The method according to Claim 23, wherein the DNS server transmits to the requesting client a redirection packet to cause the selected server to modify its respective service metric.
25. The method according to Claim 19, wherein the candidate server list includes extra, randomly selected servers selected from among the servers represented by the weights beyond the number servers addresses selected based on the weights.
 - 15
26. The method according to Claim 25, wherein the extra, randomly selected servers are a fixed percentage beyond the number of server addresses selected based on the weights.
- 20 27. The method according to Claim 25, wherein the extra, randomly selected server addresses are a fixed number beyond the number of servers selected based on the weights.

28. The method according to Claim 19, wherein each server represented in the candidate server selection list is unique from each other server represented in the list.
29. The method according to Claim 19 wherein the feedback occurs according to at least one of the following: number of times selected, duration from last feedback, time of day, or requested event.
30. The method according to Claim 19 wherein the weights are based on bias factors to reduce convergence time, the bias factors including at least one of: number of times selected, moving average based on length of recording time, historical count information, time of day, , time of year, calendar event, or part of country.
31. The method according to Claim 19 wherein the probabilities sum to one.
32. The method according to Claim 19 wherein the central server includes unique vectors of server selection weights for subsets of clients.
33. The method according to Claim 19 wherein the central server includes multiple central servers organized as a distributed system.
34. The method according to Claim 19 wherein the step of selecting the preferred server comprises the step of interrogating the servers to measure at least one of the following: network performance between client and server, server congestion, or server load.
35. By a central server in a packet communication network, a method for providing a client with a list of possible optimal servers from among a plurality of servers also on the network, comprising the steps of:
maintaining weights corresponding to a plurality of servers;

in response to receiving a client request from a client, selecting a candidate server list from among the servers represented by the weights; providing the candidate server list for the requesting client; and receiving feedback related to service by the servers for maintaining the
5 weights.

36. The method according to Claim 35 further including:
 - establishing a relationship with a DNS server to have the DNS server pass to the central server requests from clients for a server known by the central server.
- 10 37. The method according to Claim 35 wherein the candidate server list includes extra, randomly selected servers selected from among the servers represented by the weights beyond the randomly selected servers selected based on the weights.
38. The method according to Claim 35 wherein the weights are based on bias factors to reduce convergence time, the bias factors including at least one of:
 - 15 number of times selected, moving average based on length of recording time, historical count information, time of day, time of year, calendar event, or geographical location.
39. The method according to Claim 35 wherein the central server includes unique vectors of weights for subsets of clients.
- 20 40. A computer program product comprising:
 - a computer readable medium for storing data; and
 - a set of computer program instructions embodied on the computer readable medium, including instructions to:
 - maintain weights related to service provided by servers;

in response to receiving a request from a client, select a candidate server list from among the servers represented by the weights based on the weights;

5 provide the candidate server selection list for the requesting client; and
receive feedback from servers to maintain the weights.

41. The computer program product according to Claim 40 further including instructions to:
 - 10 establish a relationship with a DNS server to have the DNS server pass to the central server client requests for a server known by the central server.
42. The computer program product according to Claim 41 wherein the candidate server selection list includes extra, randomly selected servers selected from among the servers represented by the weights beyond the servers selected based on the weights.
 - 15
43. The computer program product according to Claim 40 wherein the probabilities are based on bias factors to reduce conversion time, including at least one of:
 - 20 number of times selected, moving average based on length of recording time, historical count information, time of day, time of year, calendar event, or geographical location.
44. The computer program product according to Claim 40 wherein the central server includes unique vectors of weights for subsets of clients.
45. An apparatus for providing a client with a list of possible optimal servers from among a plurality of servers also on the network, the apparatus comprising:
 - 25 means for maintaining weights based on service by servers;
means for receiving a request from a client;

means for selecting a candidate server selection list based on the weights from among the servers represented by the weights; and

means for providing the candidate server selection list for the requesting client.

- 5 46. The apparatus according to Claim 45 further comprising means for establishing a relationship with a DNS server to have the DNS server pass requests from clients for a server to the means for receiving a request from a client.
47. An apparatus for providing a client with a list of possible optimal servers from among a plurality of servers also on the network, the apparatus comprising:
- 10 a processor coupled to memory storing weights related to service provided by a plurality of servers, said processor executing a computer program to:
- maintain the weights;
- receive a request from a client;
- 15 in response to the request, select a candidate server list from among the servers represented by the weights; and
- an interface coupled to the processor and the network to provide the candidate server list for the requesting client.
48. The apparatus according to Claim 47 wherein the processor establishes a relationship with a DNS server to have the DNS server pass to the processor requests from clients for a server known by the processor.
- 20
49. The apparatus according to Claim 47 wherein the candidate server list includes extra, randomly selected servers selected from among the servers represented by the weights beyond the servers selected in part based on the weights.

50. The apparatus according to Claim 47 wherein the probabilities are based on bias factors to reduce convergence time, including at least one of:

number of times selected, moving average based on length of recording time, historical count information, time of day, time of year, calendar event, or geographical location.

5 51. The apparatus according to Claim 47 wherein the memory includes unique vectors of weights for subsets of clients.

10 52. By a server in a packet communication network, a method for participating in optimizing server selection for clients from among a plurality of servers, the method comprising:

responding to a probe from a client; and
if selected to provide service as a result of responding to the probe, (i) providing requested service to the client and (ii) reporting to a central server that service is being provided to the client.

15 53. The method according to Claim 52, further including maintaining a count of the number of clients for which the server is providing service.

54. The method according to Claim 53, wherein reporting to the central server includes reporting the count to the central server.

20 55. A computer program product comprising:

a computer readable medium for storing data; and
a set of computer program instructions embodied on the computer readable medium, including instructions to:
respond to a probe from a client; and

if selected to provide service as a result of responding to the probe, (i) provide requested service to the client and (ii) report to a central server that service is being provided to the client.

- 5 56. The computer program product according to Claim 55, further including instructions to maintain a count of the number of clients for which the server is providing service.
- 10 57. The computer program product according to Claim 56, wherein the instructions to report to the central server that service is being provided includes reporting the count to the central server.
- 15 58. An apparatus for participating in optimizing server selection for clients from among a plurality of servers, the apparatus comprising:
means for responding to a probe from a client; and
if selected to provide service as a result of responding to the probe, (i)
means for providing requested service to the client and (ii) means for reporting to a central server that service is being provided to the client.
- 20 59. The apparatus according to Claim 58, further comprising means for maintaining a count of the number of clients for which the server is providing service.
60. The apparatus according to Claim 59, wherein the means for reporting includes means for reporting the count to the central server.